

PERSONALIZING UNIVERSAL RECURRENT NEURAL NETWORK LANGUAGE MODEL WITH USER CHARACTERISTIC FEATURES BY SOCIAL NETWORK CROWDSOURCING

Bo-Hsiang Tseng, Hung-yi Lee, and Lin-Shan Lee

National Taiwan University
r02942037@ntu.edu.tw, lslee@gate.sinica.edu.tw

ABSTRACT

With the popularity of mobile devices, personalized speech recognizer becomes more realizable today and highly attractive. Each mobile device is primarily used by a single user, so it's possible to have a personalized recognizer well matching to the characteristics of individual user. Although acoustic model personalization has been investigated for decades, much less work have been reported on personalizing language model, probably because of the difficulties in collecting enough personalized corpora. Previous work used the corpora collected from social networks to solve the problem, but constructing a personalized model for each user is troublesome. In this paper, we propose a universal recurrent neural network language model with user characteristic features, so all users share the same model, except each with different user characteristic features. These user characteristic features can be obtained by crowdsourcing over social networks, which include huge quantity of texts posted by users with known friend relationships, who may share some subject topics and wording patterns. The preliminary experiments on Facebook corpus showed that this proposed approach not only drastically reduced the model perplexity, but offered very good improvement in recognition accuracy in n-best rescoring tests. This approach also mitigated the data sparseness problem for personalized language models.

Index Terms— Recurrent Neural Network, Personalized Language Modeling, Social Network, LM adaptation

1. INTRODUCTION

The personalization of various applications and services for each individual user has been a major trend. Good examples include personalized web search [1, 2] and personalized recommendation systems [3, 4, 5, 6]. In the area of speech recognition, the popularity of mobile devices such as smart phones and wearable clients makes personalized recognizers much more realizable and highly attractive. Each mobile device is used primarily by a single user, and can be connected to a personalized recognizer stored in the cloud with much better performance, because this recognizer can be well-matched to the linguistic characteristics of the individual user.

In acoustic model adaptation [7, 8, 9], personalization has

been investigated for decades and has yielded very impressive improvements with many approaches based on either HMM/GMM or CD-DNN-HMM [10]. However, there has been much less work reported on language model (LM) personalization. LM adaptation has been studied for decades [11, 12, 13], but the previous works [14, 15, 16, 17, 18] primarily focused on the problem of cross-domain or cross-genre linguistic mismatch, while the cross-individual linguistic mismatch is often ignored. One good reason for this is perhaps the difficulty in collecting personalized corpora for personalized LMs. However, this situation has changed in recent years. Nowadays, many individuals post large quantities of texts over social networks, which yield huge quantities of posted texts with known authors and given friend relationships among the authors. It is therefore possible to train personalized LMs because of the reasonable assumption that users with close friend relationships may share common subject topics, wording habits, and linguistic patterns.

Personalized LMs are useful in many aspects [19, 20, 21]. In the area of speech recognition, personalization of LMs has been proposed and investigated for both N-gram-based LMs [22] and recurrent neural networks (RNNLMs) [23] in the very limited previous works. In these previous works, text posted by many individual users and other information (such as friend relationships among users) were collected from social networks. A background LM (either N-gram-based or RNN-based) was then adapted toward an individual user's wording patterns by incorporating social texts that the target user and other users had posted, considering different aspects of their relationships and similarities between the users. In these previous works, personalization was realized by training an LM for each individual. There are inevitable shortcomings with this framework. First, even with help of the social networks, it is not easy to obtain text corpora that are helpful for a particular user for adapting a background LM towards a personalized LM. As a result, the personalized LM thus obtained easily overfits to the limited data, and therefore yields relatively poor performance on the new data of the target user. Second, to train and store a personalized LM for every user is in any case time-consuming and memory-intensive, especially considering that the number of users will only increase in the future.

Considering the above-mentioned defects in the previous framework, in this paper we propose a new RNNLM-based paradigm for personalizing LMs. In conventional RNNLM [24, 25, 26], the 1-of-N encoding of each word is taken as the input of the RNN, and then given the history word sequence, RNN outputs the estimated probability distribution for the next word. In the new paradigm proposed here, however, each user is represented by a feature vector encoding some characteristics of the user, and this feature vector augments the 1-of-N encoding feature of each word. A universal RNNLM is thus trained based on the data of these user features, together with the texts over social networks by a large number of users. The standard training method is used, except now the same words produced by different users in the training set are augmented by different user characteristic features. For each new user, his characteristic feature is extracted to extend the 1-of-N word encoding, with which the universal RNNLM can be used to recognize his speech. Because the same words produced by different users are augmented with different features, given the same history word sequence, the universal RNNLM can predict different distributions of the next word for different users. In this way, the personalization can be achieved even though all users share the same universal RNNLM. This universal RNNLM trained from the social text produced by many users is less liable to overfitting because a very large training set can be obtained by aggregating the social texts of many users. Moreover, since the recognizer for each user only requires the user's characteristic features rather than an entirely new model, the new paradigm saves time during training and memory in real-world implementations. This concept of input features for personalization is similar to the i-vectors used in deep neural network (DNN) based acoustic models [27, 28], in which the i-vector of each speaker is used to extend acoustic features such as MFCC. Preliminary experiments show that the proposed method not only reduces model perplexity but also reduces word error rates in n-best rescoring tests. In addition, we find that this approach mitigates the overfitting problem for limited personalized data. can be helpful in extracting the target user's characteristic features.

2. LM PERSONALIZATION SCENARIO

Crowdsourcing [29, 30] has varying definitions and has been applied to a wide variety of tasks. For example, a crowdsourcing approach was proposed to collect queries for information retrieval considering temporal information [31]. The MIT movie browser [32, 33] build a crowd-supervised spoken language system. In this work, a cloud-based application was implemented offering users to access to their social network via voice, and was treated as a crowdsourcing platform for collecting personal data. When the user logs into his Facebook account, he may choose to grant this application the authority to collect his acoustic and linguistic data for use in personalizing the voice access service. Users who do so

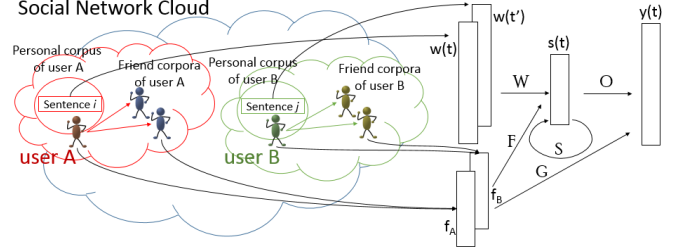


Fig. 1: The scenario for the proposed approach. When training with sentence i from user A, the user feature fed into the RNNLM can be produced by either topic distribution of user A's personal corpus or searching over user A's personal / friends corpora for sentences with topic distributions closest to this sentence i .

may enjoy the benefits of the superior recognition accuracy yielded by the personalized recognizer via the crawled data.

Fig. 1 depicts the scenario of the proposed approach. For user A, the red figure in the left part of the figure, the texts of his social network posts are crawled to form the *personal corpus* of the user (the red circle). Besides, the posts of all user A's friends (blue figures) in social network are collected to form user A's *friends corpora* (the red cloud surrounding the circle). In previous work [23] the user's personal corpus and friends corpora was used for adapting a background LM trained with a large background corpus. However, such LM adaptation suffers from overfitting due to the limited adaptation data, and as mentioned above, also incurs heavy training/memory burdens.

In this paper, instead of building a personalized RNNLM for each user, a single universal RNNLM is used by all users. As shown in the right part of Fig. 1, a corpus of posts from a large group of users serves as the training data for the universal RNNLM. This universal RNNLM comprises three layers: the input layer, the hidden layer, and the output layer, very similar to those used previously [24], except the input layer is not only the word vector $w(t)$ representing the t -th word in a sentence using an 1-of-N encoding, but concatenated with the additional user characteristic feature f . This user characteristic feature is connected to both the hidden layer $s(t)$ and output layer $y(t)$ ¹. This feature f enables the model to take into account each specific user. The network weights to be learned are the matrices W, F, S, G and O in the right part of the figure.

3. EXTRACTION OF USER CHARACTERISTIC FEATURES

We proposed two approaches to extract the user characteristic feature for each sentence, which are respectively described in Subsections 3.1 and 3.2.

¹This structure is parallel to the context dependent RNNLM variant [26] except that the context feature in the input layer is replaced by the user characteristic feature f .

3.1. User-dependent Feature

In this approach, the personal corpus for each target user is viewed as a single document, and then a topic modeling approach is used to derive the topic distribution of that document. The topic distribution of the personal corpus thus represents the language characteristics of the user and is considered as the user characteristic feature f of the user. That is, during training the universal RNNLM, the 1-of-N encoding of the words in a personal corpus are all concatenated with the same topic distribution of that personal corpus. The topic model used here is Latent Dirichlet Allocation (LDA) [34] model trained from a large corpus for many users.

3.2. Sentence-dependent Feature

Considering the fact that the personalized corpus of a user may cover many different topics, and the topic of the user may be switched dynamically and freely from one to another in the personal corpus, the topic distribution for the whole personal corpus may not very well represent each individual sentence within the personal corpus. On the other hand, even though the topic can be switched freely in the personal corpus of a user, we observe that it usually needs at least a few sentences to finish a specific topic. Therefore, to form a feature not only reflecting the characteristics of user but also a specific sentence, we can exploit a part of the personal corpus whose topic distribution is close to the sentence. This may solve the problem of mismatch between the topic distribution of the whole personal corpus and each individual training sentence.

With the above consideration, in the second approach, every sentence in the personal corpus of a user has its unique feature f which is related to not only the user but the sentence itself. In other words, the topic model is first used to infer the topic distribution of a sentence, we then use this topic distribution to search over the personal corpus of the user to find other N sentences whose topic distributions are most closest to one formed for the sentence being considered. This search process is fast since it is limited to personal corpus of the considered user only. During training the universal RNNLM, the average of the topic distributions of these N found sentences is taken as the user characteristic feature f , to be concatenated with the 1-of-N encoding features of the words in the sentence. Therefore, the same words in different sentences of a personal corpus may have different user characteristic features. We can also extend the search space to be over the friends corpora of the user as well.

The major difference of the two approaches in Subsections 3.1 and 3.2 lies in the concept of how a better language model can be obtained. In the first approach, we assume the personal corpus of a user can reflect his language characteristics, so the data for inferring the topic distribution is the whole personal corpus. In the second approach, we assume a user actually switches his topic freely from sentence to sentence, so we try to find the similar sentences to construct the user character feature to reflect the language characteristics

not only for the user but for the specific sentence itself. So, the data to form the user characteristics is limited to the N sentences found in the search process. During testing, the user characteristic feature is obtained in exactly the same way, except the N -best list of an utterance was used with the LDA model to generate the topic distribution for an utterance.

4. EFFECTS OF THE USER CHARACTERISTIC FEATURE ON RNNLM

Here we use a real example from the Facebook data to show the effect of the user characteristic features on RNNLM. User A left many posts about coffee in the Facebook data, while user B never did so. This yielded very different user characteristic features for the two users. Here the user characteristic features mentioned in subsection 3.2 by searching for the N closest sentences are used. Given the sentence “*A bottle of milk can make 3 cups of latte*” which was more likely to be produced by user A, we list in Table 1 the perplexities evaluated by a conventional RNNLM and the personalized RNNLM with different user characteristic features. The conventional RNNLM is in row (a). We see that the personalized RNNLM with the user characteristic feature f_A of user A produced a drastically decreased perplexity (152 vs 355, row (b)) because of the well-matched characteristics, while that with the user characteristic feature f_B of user B yielded a significantly increased perplexity (604 vs 355, row (c)).

Language Models	Perplexity
(a) RNNLM (conventional)	355
(b) RNNLM (with f_A)	152
(c) RNNLM (with f_B)	604

Table 1: The perplexity for sentence “*A bottle of milk can make 3 cups of latte*” using different models, where user A’s personal corpus included many posts about coffee but user B’s personal corpus contained no such posts.

5. EXPERIMENTAL SETUP

5.1. Corpus & LMs

Our experiments were conducted on a crawled Facebook corpus. A total of 42 users logged in and authorized this project to collect for research their posts and basic information. These 42 users were our target users, and were divided into 3 groups for cross validation, i.e., to train the universal LM using the data of two groups and test those for the rest. Furthermore, with their consent, the observable public data (the personal and friends corpora) of these 42 target users were also available for the experiments. This resulted in the personal data of 93,000 anonymous people and a total of 2.4 million sentences. The number of sentences for each user among the 93,000 ranged from 1 to 8,566 with a mean of 25.7, comprising 10.6 words (Chinese, English, or mixed) per sentence on average. A total of 12,000 sentences for the

42 target users was taken as the testing set, and among them 948 produced by the respective target users were taken as testing utterances for ASR experiments.

For the background corpus, 500k sentences were collected from the popular social networking site Plurk to train the topic model. Using the Mallet toolkit [35], we trained a latent Dirichlet allocation-based (LDA) topic model, taking each sentence as a document. The modified Kneser-Ney algorithm [36] was used for the N-gram LM smoothing. From the corpus the most frequent 18,000 English words and 46,000 Chinese words were selected to form the lexicon. The SRILM [37] toolkit was used for the N-gram LM training and adaptation, while RNNLM toolkit [38] was used for RNNLM here.

5.2. N-best rescoring

To generate the 1,000-best lists for rescoring, we used lattices produced using the HTK toolkit [39]. To generate the lattices we used a trigram LM adapted to the personal and friend corpora using Kneser-Ney smoothing (KN3). For first-pass decoding we used Mandarin triphone models trained on the ASTMIC corpus and the English triphone models trained on the Sinica Taiwan English corpus [40]; both corpora include hundreds of speakers. Both sets of models were adapted using unsupervised MLLR.

6. EXPERIMENTAL RESULTS

6.1. Extraction of user characteristic features

As mentioned in Section 3.2, only those N sentences most close to the sentence under consideration were used to build the user characteristic feature. Fig. 2 are the perplexities for different N (out of the user plus friend corpora) and different number of topics for LDA². The figure shows that there was almost no difference between $N = 1$ and $N = 2$, but as N increased beyond 2 the perplexity also increased, suggesting a wide variety of topics even for the same user and his friends. We thus chose $N = 1$ for the following experiments.

6.2. Perplexity

Table 2 shows the results of perplexity (PPL). Personalized Kneser-Ney tri-gram is reported in section (a) [23], where ‘B’, ‘B+P’, ‘B+P+F’ indicate respectively background (B), background plus personal corpus (B+P), and plus friends corpora in addition (B+P+F). Row (b) is RNNLM using only background corpus (‘B’) without any personalization with hidden layer size of 50 and 200. Personalizing RNNLM based on model adaption (model) [23] and the user characteristic feature (UCF) approach proposed here are respectively labeled with ‘RNN/model’ in section (c) and ‘RNN/UCF’ in section

²Only one-tenth of the personal and friend corpus was used in these preliminary experiments.

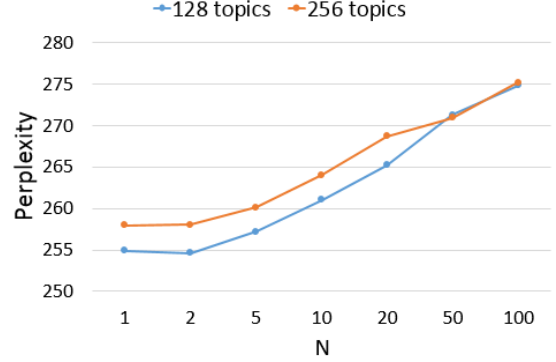


Fig. 2: Perplexities for different number of LDA topics and different number of similar sentences (N) selected to build the user characteristic feature in 3.2.

(d). In section (d), notations ‘UD’ and ‘SD’ respectively indicate extracting user-dependent (subsection 3.1) and sentence-dependent (subsection 3.2) features in the proposed approach.

We found that sentence-dependent feature outperformed the user-dependent feature ((d-1) v.s. (d-1)), which implies that the considerations of topic switching in subsection 3.2 is reasonable. Under the condition involving personal corpora (B+P), no matter using user-dependent or sentence-dependent features, the approach proposed here is always better than the model adaption approach ((d-1) (d-2) v.s. (c-1)). With the sentence-dependent features, PPL improvement is up to 102 (218 in (d-2) v.s. 320 in (c-1) in h200). This means extracting a good feature to characterize the user is more efficient than using personal data to learn a personalized RNNLM. With the friends corpora involved (B+P+F)³, the proposed approach is still better than the model adaption approach (211 in (d-3) v.s. 265 in (c-2)). When using sentence-dependent feature we may further average the found user characteristic feature with the topic distribution of the sentence being considered (RNN/UCF, SD, B+P+F, avg in (d-4)). PPL in this case can be further improved (165 in (d-4) v.s. 211 in (d-3)). With this best model obtained here (RNN/UCF, SD, B+P+F, avg), the perplexity is reduced by 58.5% compared to RNNLM without personalization (RNN, B in (b)), 37.7% compared to the model adaption approach with friends corpora used (RNN/model, B+P+F in (c-2)).

6.3. Word error rate (WER)

Table 3 reports the word error rates (WER) with the same notation as in Table 2. Section (a) is for the three different tri-gram LMs without and with personalization. As expected, with more adaptation data, the tri-gram LMs performed better ((a-3) < (a-2) < (a-1)) [23]. We used the best adapted tri-gram LM (KN3, B+S+F in (a-3)) to generate 1000-best lists for RNNLM rescoring. Section (b) is for rescoring results using RNNLM without personalization, while sections (c) and (d)

³When extracting the sentence-dependent feature, search space is over both personal and friends corpora.

Perplexity		h50	h200
(a)	(a-1) KN3, B	343	
	(a-2) KN3, B+P	299	
	(a-3) KN3, B+P+F	233	
(b)	RNN, B	441	398
(c)	(c-1) RNN/model, B+P	350	320
	(c-2) RNN/model, B+P+F	296	265
(d)	(d-1) RNN/UCF, UD, B+P	313	270
	(d-2) RNN/UCF, SD, B+P	269	218
	(d-3) RNN/UCF, SD, B+P+F	229	211
	(d-4) RNN/UCF, SD, B+P+F, avg	192	165

Table 2: Perplexity (PPL) Results. KN3 represents Kneser-Ney tri-gram, while ‘RNN/model’ and ‘RNN/UCF’ are for Personalizing RNNLM based on model adaption (model) and the user characteristic feature (UCF) approach proposed here respectively. Notation ‘B’, ‘B+P’ and ‘B+P+F’ respectively indicate using only background corpus (B), plus personal corpus (B+P), and plus friends corpora in addition (B+P+F). Notation ‘UD’ and ‘SD’ respectively indicate extracting user-dependent and sentence-dependent features in the proposed approach. The results for RNNLM with hidden layer size of 50 and 200 are listed.

are for model adaption (model) approach and the proposed approach (UCF) respectively. For sentence-dependent (SD) features, we viewed the 1000-best list as a single document and used the LDA topic model to infer the topic distribution, and then search for the closest sentences as mentioned in subsection 3.2 to construct the sentence-dependent feature of each utterance for rescoring. For user-dependent (UD) features in the proposed approach, because the feature is extracted from the personal corpus of the user and independent of the input utterance, so the feature extraction process does not depend on ASR. Regardless of the features used and the data involved, the proposed approach was always better than the model adaption approach ((d-1, 2) v.s. (c-1) and (d-3, 4) v.s. (c-2)).

To our surprise, for 200 hidden layer units the proposed approach with user-dependent feature is better than sentence-dependent feature in terms of WER ((d-1) v.s. (d-2, 3, 4) for h200). This may be because the user-dependent feature is estimated from the training corpus of target user thus not influenced by ASR errors at all; while for sentence-dependent feature, the topic distribution from N-best list was inaccurate due to ASR errors. This was verified in the oracle experiments in section (e), in which we used the topic distribution of the reference transcription of the utterance to replace the topic distribution of N-best list and do the rescoring⁴. Here we see the sentence-dependent feature (SD) is better than user-dependent feature (UD) ((e-2) v.s. (e-1) while (e-3) (e-4) are even better). Also, with topic distributions from the refer-

⁴In the oracle experiments in section (e), results of user-dependent feature (e-1) were the same as those in (d-1) because ASR was not involved in extracting the feature.

WER (%)		h50	h200
(a)	KN3, B	43.80	
	(a-1) KN3, B+P	43.39	
	(a-2) KN3, B+P+F	41.95	
(b)	(a-3) RNN, B	41.12	41.14
(c)	(c-1) RNN/model, B+P	40.84	40.87
	(c-2) RNN/model, B+P+F	40.71	40.68
(d)	(d-1) RNN/UCF, UD, B+P	40.48	40.16
	(d-2) RNN/UCF, SD, B+P	40.47	40.36
	(d-3) RNN/UCF, SD, B+P+F	40.43	40.40
	(d-4) RNN/UCF, SD, B+P+F, avg	40.23	40.26
oracle	(e-1) RNN/UCF, UD, B+P	40.48	40.16
	(e-2) RNN/UCF, SD, B+P	40.15	40.09
	(e-3) RNN/UCF, SD, B+P+F	40.03	39.95
	(e-4) RNN/UCF, SD, B+P+F, avg	39.40	39.45

Table 3: Word error rate (WER) results with same notations as in Table 2. For sentence-dependent (SD) features, the topic distributions are estimated from N-best lists in section (d), while from reference transcriptions in section (e) (oracle).

ence transcriptions, the results of sentence-dependent feature can be improved by absolute 0.81% (from 40.26% in (d-4) to 39.45% in (e-4) for h200).

So for the real best result here (RNN/UCF, UD, B+P for h200 in (d-1)), we reduced WER by 1.79% compared to the best of KN3 including friends corpora (41.95% in (a-3)), from which the 1000-best lists were obtained, 0.98% compared to RNNLM without personalization (41.14% in (b)), 0.52% compared to the best of model adaption approach (40.68% in c-2). In the oracle case the best result can be even much better (39.45% RNN/UCF, SD, B+P+F, avg in (e-4)), which indicates the space for further improvement.

6.4. Analysis

6.4.1. WER over all target users

Because the average didn’t tell whether the proposed approach is actually helpful for most users or for just very few users, we plot in addition the WER change obtained across the all 42 target users in Fig. 3. The three figures in the upper row compare respectively the proposed approach with user-dependent feature (RNN/UCF, UD, B+P, in (d-1) of Table 3), sentence-dependent feature (RNN/UCF, SD, B+P+F, avg, in (d-4) of Table 3) and sentence-dependent feature in oracle experiments (RNN/UCF, SD, B+P+F, avg, in (e-4)) with the baseline of RNNLM without personalization (row (b) in Table 3), and the three figures in lower row compare the same three approaches with the model adaption approach (row (c-1) in Table 3). Each figure has 42 bars for the 42 target users, sorted based on the WER change. Here a negative value means that the proposed approach here offered WER reduction to the user. From Fig. 3, we see the proposed approach offered better performance to most target users. For example, in the first figure on the upper row (

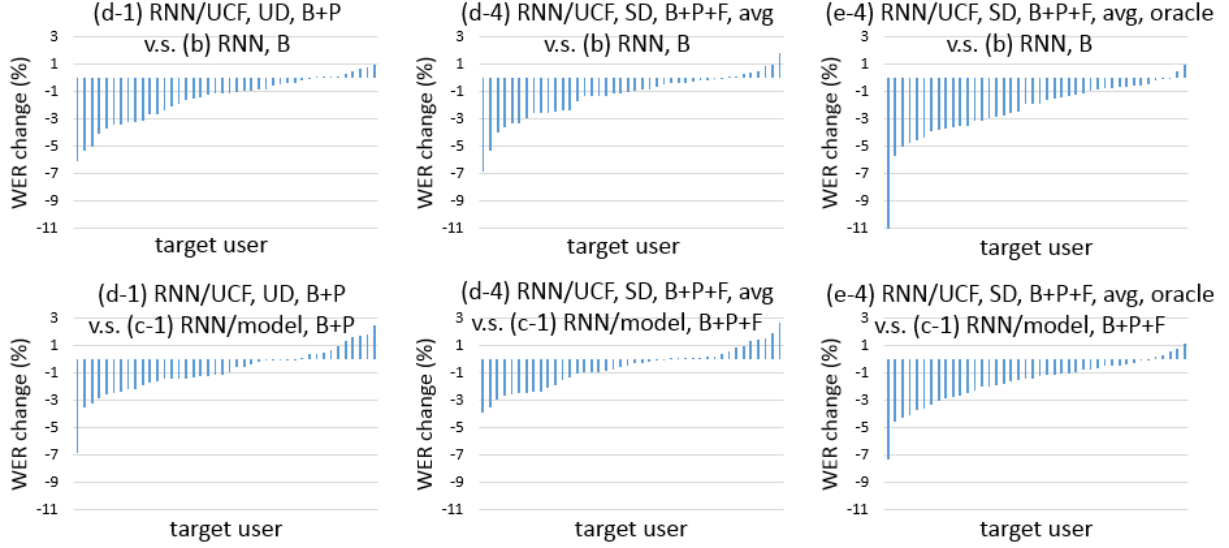


Fig. 3: WER changes across all 42 target users. The three figures in the upper row compare respectively the proposed approach with user-dependent feature (RNN/UCF, UD, B+P, in (d-1) of Table 3), sentence-dependent feature (RNN/UCF, SD, B+P+F, avg, in (d-4) of Table 3) and sentence-dependent feature in oracle experiments (RNN/UCF, SD, B+P+F, avg, in (e-4)) with the baseline of RNNLM without personalization (row (b) in Table 3), and the three figures in lower row compare the same three approaches with the model adaption approach (row (c-1) in Table 3).

(d-1) RNN/UCF, UD, B+P v.s. (b) RNN,B), 9 users had worse WER with our approach, all by less than 1%, but all other users had WER reduction, 24 of them by more than 1%. Similar for the rest cases. The results show the proposed approach offered improvements to most target users.

6.4.2. Size of personal corpus

As mentioned above, the model adaptation approach results in overfitting to the limited personal data and may yield poor performance on a particular user’s new data. This is illustrated in Fig. 4. The horizontal axis of the figure is the percentage of the original personal corpus used, where 1.00 means using the entire original personal corpus, that is, those cases (c-1) RNN/model, (d-1) RNN/UCF, UD, B+P and (d-2) RNN/UCF, SD, B+P in Tables 2 and 3 for h50. We see that as less data were available, the proposed approach (d-1) and (d-2) demonstrated much smaller increases in perplexity and much more stable WER, whereas for the model adaptation approach (c-1), the perplexity and WER increased significantly at a greater rate. The result of different size of friends corpora has the same trend.

7. CONCLUSIONS

In this paper, we proposed a new framework for personalizing a universal RNNLM using data crawled over social networks. The proposed approach is based on a user characteristic feature extracted from the user corpus and friends corpora, which is not only user-dependent but sentence-dependent fea-

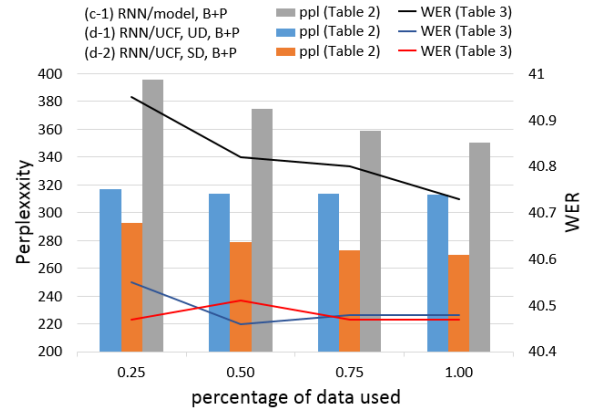


Fig. 4: Perplexity (PPL) and word error rate (WER) for different sizes of personal corpus for the model adaptation approach ((c-1) RNN/model, B+P) and the proposed approaches ((d-1) RNN/UCF, UD, B+P and (d-2) RNN/UCF, SD, B+P). The horizontal axis is the percentage of the original personal corpora used (1.00 on the right is for the data in Tables 2 and 3).

ture. This universal RNNLM can predict different word distributions for different users given the same context. Experiments demonstrated really good improvements in both perplexity and WER, and the proposed approach is much more robust to data sparseness than the previous work.

8. REFERENCES

- [1] Geoffrey Zweig and Shuangyu Chang, "Personalizing model m for voice-search.," in *INTERSPEECH*, 2011, pp. 609–612.
- [2] Micro Speretta and Susan Gauch, "Personalized search based on user search histories," in *Web Intelligence, 2005. Proceedings. The 2005 IEEE/WIC/ACM International Conference on*. IEEE, 2005, pp. 622–628.
- [3] Yoon Ho Cho, Jae Kyeong Kim, and Soung Hie Kim, "A personalized recommender system based on web usage mining and decision tree induction," *Expert systems with Applications*, vol. 23, no. 3, pp. 329–342, 2002.
- [4] Yehuda Koren, Robert Bell, and Chris Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, , no. 8, pp. 30–37, 2009.
- [5] Frank Edward Walter, Stefano Battiston, and Frank Schweitzer, "A model of a trust-based recommendation system on a social network," *Autonomous Agents and Multi-Agent Systems*, vol. 16, no. 1, pp. 57–74, 2008.
- [6] Moon-Hee Park, Jin-Hyuk Hong, and Sung-Bae Cho, "Location-based recommendation system using bayesian users preference model in mobile devices," in *Ubiquitous Intelligence and Computing*, pp. 1130–1139. Springer, 2007.
- [7] Christopher J Leggetter and Philip C Woodland, "Maximum likelihood linear regression for speaker adaptation of continuous density hidden markov models," *Computer Speech & Language*, vol. 9, no. 2, pp. 171–185, 1995.
- [8] Phil C Woodland, "Speaker adaptation for continuous density hmms: A review," in *ISCA Tutorial and Research Workshop (ITRW) on Adaptation Methods for Speech Recognition*, 2001.
- [9] Jean-Luc Gauvain and Chin-Hui Lee, "Maximum a posteriori estimation for multivariate gaussian mixture observations of markov chains," *Speech and audio processing, iee transactions on*, vol. 2, no. 2, pp. 291–298, 1994.
- [10] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al., "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *Signal Processing Magazine, IEEE*, vol. 29, no. 6, pp. 82–97, 2012.
- [11] Roland Kuhn and Renato De Mori, "A cache-based natural language model for speech recognition," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 12, no. 6, pp. 570–583, 1990.
- [12] Rukmini M Iyer and Mari Ostendorf, "Modeling long distance dependence in language: Topic mixtures versus dynamic cache models," *Speech and Audio Processing, IEEE Transactions on*, vol. 7, no. 1, pp. 30–39, 1999.
- [13] Ronald Rosenfeld, "A maximum entropy approach to adaptive statistical language modelling," *Computer Speech & Language*, vol. 10, no. 3, pp. 187–228, 1996.
- [14] Jerome R Bellegarda, "Statistical language model adaptation: review and perspectives," *Speech communication*, vol. 42, no. 1, pp. 93–108, 2004.
- [15] Aaron Heideel and Lin-shan Lee, "Robust topic inference for latent semantic language model adaptation," in *Automatic Speech Recognition & Understanding, 2007. ASRU. IEEE Workshop on*. IEEE, 2007, pp. 177–182.
- [16] Bo-June Paul Hsu and James Glass, "Style & topic language model adaptation using hmm-lda," in *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2006, pp. 373–381.
- [17] David Mrva and Philip C Woodland, "Unsupervised language model adaptation for mandarin broadcast conversation transcription.," in *INTERSPEECH*. Citeseer, 2006.
- [18] Yik-Cheung Tam and Tanja Schultz, "Correlated latent semantic model for unsupervised lm adaptation," in *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*. IEEE, 2007, vol. 4, pp. IV–41.
- [19] Yu-Yang Huang, Rui Yan, Tsung-Ting Kuo, and Shou-De Lin, "Enriching cold start personalized language model using social network information," *ACL14*, pp. 611–617, 2014.
- [20] Gui-Rong Xue, Jie Han, Yong Yu, and Qiang Yang, "User language model for collaborative personalized search," *ACM Transactions on Information Systems (TOIS)*, vol. 27, no. 2, pp. 11, 2009.
- [21] Arjumand Younus, Colm ORiordan, and Gabriella Pasi, "A language modeling approach to personalized search based on users microblog behavior," in *Advances in Information Retrieval*, pp. 727–732. Springer, 2014.
- [22] Tsung-Hsien Wen, Hung-Yi Lee, Tai-Yuan Chen, and Lin-Shan Lee, "Personalized language modeling by crowd sourcing with social network data for voice access of cloud applications," in *Spoken Language Technology Workshop (SLT), 2012 IEEE*. IEEE, 2012, pp. 188–193.

- [23] Tsung-Hsien Wen, Aaron HeideI, Hung-yi Lee, Yu Tsao, and Lin-Shan Lee, "Recurrent neural network based language model personalization by social network crowdsourcing.," in *INTERSPEECH*, 2013, pp. 2703–2707.
- [24] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur, "Recurrent neural network based language model.," in *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010*, 2010, pp. 1045–1048.
- [25] Tomáš Mikolov, Stefan Kombrink, Lukáš Burget, Jan Honza Černocký, and Sanjeev Khudanpur, "Extensions of recurrent neural network language model," in *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*. IEEE, 2011, pp. 5528–5531.
- [26] Tomas Mikolov and Geoffrey Zweig, "Context dependent recurrent neural network language model.," in *SLT*, 2012, pp. 234–239.
- [27] George Saon, Hagen Soltau, David Nahamoo, and Michael Picheny, "Speaker adaptation of neural network acoustic models using i-vectors," in *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*. IEEE, 2013, pp. 55–59.
- [28] Vishwa Gupta, Patrick Kenny, Pierre Ouellet, and Themis Stafylakis, "I-vector-based speaker adaptation of deep neural networks for french broadcast audio transcription," in *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. IEEE, 2014, pp. 6334–6338.
- [29] Anhai Doan, Raghu Ramakrishnan, and Alon Y Halevy, "Crowdsourcing systems on the world-wide web," *Communications of the ACM*, vol. 54, no. 4, pp. 86–96, 2011.
- [30] Robert Munro, Steven Bethard, Victor Kuperman, Vicky Tzuyin Lai, Robin Melnick, Christopher Potts, Tyler Schnoebelen, and Harry Tily, "Crowdsourcing and language studies: the new generation of linguistic data," in *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk*. Association for Computational Linguistics, 2010, pp. 122–130.
- [31] Klaus Berberich, Srikanta Bedathur, Omar Alonso, and Gerhard Weikum, *A language modeling approach for temporal information needs*, Springer, 2010.
- [32] Jingjing Liu, Scott Cyphers, Panupong Pasupat, Ian McGraw, and Jim Glass, "A conversational movie search system based on conditional random fields.," in *INTERSPEECH*, 2012.
- [33] Ian McGraw, Scott Cyphers, Panupong Pasupat, Jingjing Liu, and Jim Glass, "Automating crowd-supervised learning for spoken language systems.," in *INTERSPEECH*, 2012.
- [34] David M Blei, Andrew Y Ng, and Michael I Jordan, "Latent dirichlet allocation," *the Journal of machine Learning research*, vol. 3, pp. 993–1022, 2003.
- [35] Andrew Kachites McCallum, "Mallet: A machine learning for language toolkit," <http://mallet.cs.umass.edu>, 2002.
- [36] Frankie James, "Modified kneser-ney smoothing of n-gram models," *Research Institute for Advanced Computer Science, Tech. Rep. 00.07*, 2000.
- [37] Andreas Stolcke et al., "Srilm-an extensible language modeling toolkit.," in *INTERSPEECH*, 2002.
- [38] Tomas Mikolov, Stefan Kombrink, Anoop Deoras, Lukar Burget, and Jan Cernocký, "Rnnlm-recurrent neural network language modeling toolkit," in *Proc. of the 2011 ASRU Workshop*, 2011, pp. 196–201.
- [39] Steve Young, Gunnar Evermann, Mark Gales, Thomas Hain, Dan Kershaw, Xunying Liu, Gareth Moore, Julian Odell, Dave Ollason, Dan Povey, et al., *The HTK book*, vol. 2, Entropic Cambridge Research Laboratory Cambridge, 1997.
- [40] Ching-Feng Yeh, Aaron HeideI, Hong-Yi Lee, and Lin-Shan Lee, "Recognition of highly imbalanced code-mixed bilingual speech with frame-level language detection based on blurred posteriorgram," in *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*. IEEE, 2012, pp. 4873–4876.